

**Braakmann N, Wildman JR.**

**[Reconsidering the impact of family size on labour supply: The twin-problems of the twin-birth instrument.](#)**

***Journal of the Royal Statistical Society: Series A* 2016**

**DOI: 10.1111/rssa.12160**

**Copyright:**

This is the peer reviewed version of the above article, which has been published in final form at <http://dx.doi.org/10.1111/rssa.12160>. This article may be used for non-commercial purposes in accordance with [Wiley Terms and Conditions for Self-Archiving](#).

**Date deposited:**

12/10/2015

**Embargo release date:**

22 January 2017

# **Reconsidering the impact of family size on labour supply: The twin problems of the twin-birth instrument**

Nils Braakmann and John Wildman

Newcastle University\*

[This version: September 28, 2015]

## Abstract

We consider two econometric problems when investigating the impact of family size on labour market outcomes using the popular twin-birth instrument. The first is the potential for omitted variable bias caused by the fact that fertility treatments are linked to twin births and are typically unobserved. We present estimates corrected for this bias and find it to be comparatively small. Second, we show that the effects of twin-birth induced variation in family size, as well as characteristics of the compliers, vary substantially with time passed since birth, which has consequences for the interpretation of estimates across samples and time.

---

\* Both Newcastle University, Business School – Economics, 5 Barrack Road, Newcastle upon Tyne, NE1 4SE, UK. Email: [nils.braakmann@newcastle.ac.uk](mailto:nils.braakmann@newcastle.ac.uk); [john.wildman@newcastle.ac.uk](mailto:john.wildman@newcastle.ac.uk).

All estimates used Stata 13.1. Do files are available from the first author on request. A previous version of this paper was circulated as “Fertility treatments and the use of twin births as an instrument for family size”. We thank the associate editor, three anonymous reviewers, Marco Alfano, William H. Green, Victor Lavy, Christian Merkl, Steffen Mueller, Regina Riphahn, seminar participants at the DIW Berlin, in Hannover, Lancaster, Lueneburg, Manchester, Newcastle and Nuremberg as well as participants at the 2014 EALE and ESPE conferences and the 2015 annual meetings of the Royal Economic Society and the Scottish Economic Society for comments.

**Keywords:** Twin-birth instrument, labour supply, fertility

**JEL-classification:** C26, J13, J22

## **1 Introduction**

Estimating the impact of family size on labour supply outcomes is complicated by endogeneity problems that necessitate the use of instrumental variable (IV) methods or related strategies. A popular instrument in this literature (e.g., Bronars and Grogger, 1994; Angrist and Evans, 1998, Jacobsen et al., 1999) as well as in the literature on the link between family size and children's outcomes (e.g., Rosenzweig and Wolpin, 1980; Black, Devereux and Salvanes, 2005 and Angrist, Lavy and Schlosser, 2010) is the occurrence of twin births. The occurrence of a twin birth, on face value, looks like the perfect candidate for an instrument - it is clearly correlated with family size and it appears reasonable that it affects labour supply only through family size. However, there are two potential problems with using twin births as an instrument, both of which are related to the link between fertility treatments and multiple births documented in the medical literature. The first is the potential of omitted variable bias caused by the fact that fertility treatments are typically unobserved. We present estimates corrected for this bias and find it to be comparatively small. The second issue is more subtle. We show that the impact of a twin birth on family size (the first stage) and the impact on labour supply (the reduced form) vary substantially with time passed since the birth of the twins. We also show that the characteristics of the compliers, those individuals who end up with a larger-than-planned family size because of the twin birth, change substantially with time passed since birth. As a consequence estimates from models using the twin-birth instrument with a single cross-section of data, such as a census, depend to some extent on the age distribution of twins in the data. Twin births are not uniformly distributed across time but, since the early 1980s (see Figure 1), are increasing both in

absolute number and as the share of all maternities. This increase is likely to be partially caused by the changing prevalence of fertility treatments, making it difficult to compare results based on different samples.

(FIGURE 1 ABOUT HERE.)

In this paper we use data from the first 3 sweeps of the British Millennium Cohort Study (MCS) that follows a random sample of babies and their mothers born during late 2000 and 2001 (see section 2 for details on the data). In a first step, we consider a threat to the future, though not necessarily past, validity of the twin-birth instrument, namely the increasing use of fertility treatments, such as in-vitro fertilization (IVF) or drug treatment with Clomiphene citrate. Within the UK the use of fertility treatments has increased in most years since 1991. For IVF, in 1991 there were around 8,000 cycles, by 2011 this had increased to just over 60,000 (Human Fertilisation and Embryology Authority, 2012). In 1992, in the UK, 0.3% of all babies born resulted from IVF treatment, by 2010 this had increased to 2% (Human Fertilisation and Embryology Authority, 2012).

The problem for the twin-birth instrument arises because fertility treatments greatly increase the risk of a multiple birth, a fact well-established in the medical literature (e.g., Callahan et al., 1994; Gleicher et al., 2000; Fauser, Devroey and Macklon, 2005). In the dataset used in this paper we find that the probability of having either twins or triplets increases from around 1% for women without fertility treatment, to about 13% for women with fertility treatment. Furthermore, 24% of all the multiple births we observe in our sample are to women who have received fertility treatment, despite them comprising only 2.6% of our sample. In the following we will generally use “multiple births” or “twin births” interchangeably. “Twin birth is the expression commonly used in the economics literature. “Multiple births” would

strictly be more accurate as people might give birth to triplets, quadruplets, etc. However, this distinction makes little difference in practice as 96% of multiple births in our sample are twins.

The link between fertility treatments and twin births and the potential threat for the use of the latter as an instrumental variable has been discussed in the literature (see, e.g., Angrist, Lavy and Schlosser, 2010, p. 798, who discuss a potential bias arising from fertility treatments and then use a sample restricted to a time period before fertility treatments became common in Israel), but its actual impact has yet to be quantitatively analysed. The main theoretical concern is that while multiple births are probably still more or less random conditional on having received fertility treatment, they are unlikely to be unconditionally random.

More problematically, deciding to undergo fertility treatment is a choice that is likely to be correlated with a number of characteristics that also influence labour supply – most prominently a very strong desire for children, but, as we demonstrate later in this paper, also with factors such as age, education, having worked before pregnancy, being white, marriage, family planning, complications during the pregnancy (i.e., health) and the birth weight of the first-born/only child. Comparisons of labour supply and other characteristics in all sweeps of our data suggest mothers with and without fertility treatment are different from each other, regardless of the number of children resulting from the pregnancy. Given that we do not observe fertility treatments in most datasets commonly used by economists, these differences will introduce correlation between multiple births and (unobserved) determinants of fertility, which will render the instrument endogenous.

In our view these issues do not necessarily invalidate some of the earlier results in the literature and may not invalidate the future use of this instrument in

countries or time periods where fertility treatments are relatively uncommon. However, fertility treatments may pose a threat for the future use of this instrument in countries where they occur regularly and, more importantly, where multiple births resulting from fertility treatments are quantitatively important. Efforts to reduce the occurrence of multiple births from fertility treatments that are under way in a number of countries, including the UK, may also facilitate the future use of the twin birth-instrument.

We estimate first stages and labour supply regressions (second stages) for six models: Our base specification is one that could be estimated using most household datasets where information on fertility treatments is missing, i.e., we use the birth of twins or triplets as an instrument for family size on a range of outcomes relating to labour supply. In a second model, we additionally condition on having received fertility treatment. A comparison of these two models allows us to quantify (and correct for) the bias in the estimates in the base model. As fertility treatments are typically unobserved in most datasets, we estimate a third model that instead conditions on a set of commonly observed variables that we know to differ between women with and without fertility treatments. Results from this model allow us to make statements about whether this conditioning strategy might be a feasible approach when information on fertility treatments is lacking. In a fourth model, we condition on both fertility treatments and the same characteristics used in the previous model. This specification allows us to check whether the correlation between pre-pregnancy characteristics and multiple births arises exclusively because of fertility treatments. Finally, we investigate whether labour supply responses differ between women with and without fertility treatment by estimating separate regressions for these two groups. Our findings suggest that the instrument generally becomes

stronger in the first stages after conditioning on fertility treatments, while the second stage results across all models are qualitatively identical, i.e., the estimates always have the same sign, with only small changes in magnitude.

In a second contribution, we demonstrate that the impact of the twin-birth-induced variation in family size on labour supply depends crucially on the time passed since the occurrence of the twin births. We rely on the first three sweeps of the MCS with interviews conducted 9 months (sweep I), 3 years (sweep II) and 5 years (sweep III) after birth. We find that the impact of a twin birth on family size (the first stage) weakens over time, which is consistent with individuals adjusting their future fertility after the random shock of a multiple birth. First stages across all 3 sweeps continue to show a strong positive relationship between the occurrence of a multiple birth and family size.

As well as family size adjustments there are other factors that may contribute to changes in the results across sweeps. Our results suggest that there are major changes in the composition of the complier group, i.e., those individuals who end up with a larger-than-planned family because of a multiple birth, across the three sweeps. Furthermore, we can expect the reduced form - the impact of a twin birth on labour supply - to vary over time as the twins grow up, attend school and (at some stage) leave their parents' home. Consequently, second stages - the ratio of the reduced form and the first stage - differ substantially across the three sweeps. Specifically, there are strong negative effects of the twin-birth induced variation in family size on the mother's employment probability after 9 months. These become weaker after 3 years and essentially disappear after 5 years, which coincides with the children entering school.



These time-varying treatment effects would be comparatively innocuous if the share of twin births was constant over time. This is, however, unlikely to be the case for at least two reasons: First, there is a general trend towards giving birth later in life in many societies. As older mothers are more likely to give birth to twins or triplets (e.g., Black, Devereux and Salvanes, 2005), this is likely to make twins more common in more recent years. Second, the availability and price of fertility treatments, as well as their link to twin births (due to improved medical treatments), also differs widely across time. These two factors will lead to problems of comparability when considering labour supply estimates based on the twin-birth instrument estimated on different cross-sections. Estimates based on a single cross-section such as a census identify some weighted average of these time-varying treatment effects, where the weights depend on the age distribution of twins in the dataset that is used. As different cross-sections are likely to have different distributions of twins, it is possible that the effects will differ across datasets, even in cases where individual-level effects are identical. This in turn makes comparisons between papers using different samples complicated as it adds another source of heterogeneity.

The remainder of this paper is organised as follows: Section 2 describes the data, section 3 explains the methodology underlying our paper. Results can be found in section 4. Section 5 discusses implications for the future use of the multiple-birth instrument. Section 6 concludes.

## **2 Data**

We use data from three sweeps of the Millennium Cohort Study (MCS), which tracks a random sample of children (and their families) born in the UK in 2000-2001.

Interviews were conducted at sweep one when the children were around 9 months old, subsequent interviews took place when the children were 3 and 5 years old. Details on the design and sampling in the MCS can be found in Dex and Joshi (2005) and Hansen and Joshi (2007). The dataset is one of the few that we are aware of that includes information on fertility treatments as well as information on the mother and the development of the child. The dataset only contains mothers with at least one child, which is the group where the instrument has predictive power for family size: The multiple-birth instrument cannot be used to model the decision whether someone has one vs. no child, as everyone who gives birth to twins or triplets will have decided to have at least one child. It has predictive power for the number of children beyond one since someone who planned to have one (additional) child will end up with two or three instead.

Our estimation sample is based on the following restrictions: First, we use only cases where the mother conducted the parent interview, leading to the loss of 28 observations where the father was interviewed. Second, the MCS tracks the children born during the sampling week, not necessarily the parents, i.e., the main respondent can change in each sweep, either because the partner was interviewed or because the main carer for the child changed, for example because of adoption or death. For sweeps 2 and 3 we only use cases where the same person as in sweep 1 was interviewed, resulting in the loss of 881 (from 15,590) observations in sweep 2 and 226 (from 12,984) observations in sweep 3. We also lose some observations in each sweep due to missing values (around 150 observations each in sweeps 1 and 2 and around 100 in sweep 3). Following these restrictions we have 18,340 observations for sweep 1, 14,460 for sweep 2 and 12,581 for sweep 3. We also repeated all estimations on a balanced sample, which did not substantially change the results.

Our main outcomes of interest are i) various dummies for employment status, mainly whether the mother is working, self-employed, a student or at home to care for the family, ii) the mother's weekly working hours, calculated in two ways, either with zeros or with missing values for people not working, and finally, iii) whether she has a partner who is working. In sweep 1, we additionally have information on whether she is currently on maternity leave.

We have two variables of interest: The first is whether the mother gave birth to twins or triplets. Almost all multiple births in the dataset are twins with only 10 cases of triplets. The latter are split equally between women with and without fertility treatment. Our sample contains 254 multiple births (i.e., twins or triplets) in sweep 1, of these, 193 appear in sweep 2 and 170 appear in sweep 3. Our second key variable is whether the pregnancy was preceded by fertility treatment. In sweep 1, we have 478 women with fertility treatments, of these 394 remain in sweep 2 and 348 in sweep 3. The most common fertility treatment in the data is drug therapy with Clomiphene citrate, followed by various forms of IVF. All of these are associated with a higher frequency of multiple births relative to births not preceded by fertility treatments, but to varying degrees: The probability of a multiple birth after treatment with Clomiphene citrate is 9%, which increases to 23% after in vitro fertilization. For untreated women the corresponding probability is 1%. The variable of interest in all second stage regressions is the number of children each woman has given birth to at each sweep. Note that women can have other children than the one tracked by the MCS.

Table 1 presents descriptive information on the estimation sample.

[TABLE 1 ABOUT HERE.]

### **3 Twin births as an instrument for fertility**

### *A The basic identification strategy and the twin-birth instrument*

To illustrate the basic identification problems we use a causal diagram (or directed acyclic graph (DAG)) (Pearl, 2000; see Morgan and Winship, 2007, for a textbook treatment). In Figure 2 each directed edge (i.e., single headed arrow) such as the one from *family size* to  $Y$  represents a cause-effect relationship between variables in the model, in the sense that the variable at the origin of the edge (start of the arrow) causes the variable at the terminus. A bidirected edge, such as the one between  $X_1$  and  $X_2$ , represents common causes of the two factors that are not part of the model.

(FIGURE 2 ABOUT HERE.)

In Figure 2 we are interested in the link between *family size* and  $Y$  that can be written as a linear equation

$$Y_i = \alpha + \tau * family\ size_i + \varepsilon_i, \quad (1)$$

where  $\tau$  is the parameter of interest. In female labour supply regressions,  $Y_i$  would typically either be a dummy for labour force status or some other measure of labour supply such as desired or actual working hours, while *family size* <sub>$i$</sub>  would typically be the mother's number of children, or the number of children that live in the same household as her.

A direct estimation of this link is hindered by the presence of (potentially unobserved) confounding variables,  $X$ . Clearly, if all variables in  $X$  were observed, it would be possible to condition on them and use OLS, matching or other selection-on-observables estimators to look at the link between family size and the outcome. In the more realistic case where some variables are unobserved, these would be part of  $\varepsilon_i$  and would render *family size* <sub>$i$</sub>  endogenous. For example, in female labour supply models, both family size and the propensity to work will be influenced by (typically

unobserved) preferences for work and family size. Furthermore, a woman's work opportunities will to some extent determine the opportunity costs of childrearing.

If we ignore the issues caused by fertility treatments, one way to proceed is to use the occurrence of a *multiple birth* as an instrument for *family size*. This appears to be an attractive strategy because the biological process governing whether a pregnancy results in a singleton or a multiple birth is outside of the control of the respective parents and thus uncorrelated with any unobserved preferences for family life, any parental optimization process, or the opportunity costs of childrearing. The only further consideration required is that the age of the mother is included in the model, as it is known that multiple births become more likely for older mothers (e.g., Black, Devereux and Salvanes, 2005). It is comparatively easy to account for this by conditioning on age in a flexible way, for example through age dummies.

In a heterogeneous effects framework (Angrist and Imbens, 1994; Angrist, Imbens and Rubin, 1996) the resulting estimates are interpreted as the local average treatment effect (LATE) for those people who end up with a larger-than-planned family due to the multiple birth (the compliers). Examples would be people who planned to have two children, but then have twins at the second pregnancy or someone who wanted one child, but ended up with twins. However, not everyone who gives birth to twins will end up with a larger than planned family. For example, someone wanting two children whose first pregnancy results in twins would have matched realized and planned family size, if they have no further children. We will return to this issue in subsection 3.C below.

This general scenario is summarised in Figure 2, panel (a): A *multiple birth* leads to quasi-random variation in *family size* that is unrelated to the confounders  $X$

(or equivalently to  $\varepsilon$ ). In this case the probability limit of the IV estimate of  $\tau$  can be written as:

$$\hat{\tau} = \tau + \frac{Cov(multiple\ birth, \varepsilon)}{Cov(multiple\ birth, family\ size)} \quad (2)$$

### *B. Omitted variable bias through fertility treatments*

Equation (2) makes it clear that if multiple births and the unobservables,  $\varepsilon_i$ , from (1) are uncorrelated, the IV estimate will be consistent as  $Cov(multiple\ birth, \varepsilon)$  would be zero and the bias term in equation (2) would disappear. A central condition for this to be plausible is that twin births are (more or less) random. However, with fertility treatments this is unlikely to be the case. Fertility treatments are known to result in multiple births and fertility treatments are likely to be correlated with at least some of the confounders: In many countries, fertility treatment is expensive and not fully covered by (state) health insurance, which implies that it is likely to be correlated with parental resources. These in turn matter for labour supply and parental investment in children as they determine the budget constraint and the (non-labour) income a parent can expect when not working. Furthermore, pregnancies preceded by fertility treatment are by definition always planned. They are also likely to be correlated with a strong desire for children as fertility treatments are generally preceded by a number of attempts to conceive naturally.

Panel (b) of Figure 1 illustrates the resulting problem: Fertility treatments create an association between *multiple birth* and the confounders in  $X$ , i.e., multiple births are not randomly assigned. This in turn opens a backdoor path  $Y \leftarrow X \rightarrow fertility\ treatment \rightarrow multiple\ birth \rightarrow family\ size \rightarrow Y$  between *multiple birth* and the outcome. In more standard econometric terms, we can consider fertility treatments as an omitted variable. This means that the error term for equation (1) can be re-written as

$$\varepsilon_i = \delta_I * fertility\ treatment_i + v_i \quad (3)$$

where  $\delta_I$  is the marginal effect of fertility treatment on labour market decisions and  $v_i$  is a new error term that is still correlated with family size since it is likely that family size will still be endogenous after conditioning on having received fertility treatment. From (3) we can see that the covariance between multiple birth and  $\varepsilon_i$  is

$$Cov(multiple\ birth, \varepsilon) = \delta_I * Cov(multiple\ birth, fertility\ treatment) \quad (4)$$

Using (4) we can write the probability limit of  $\tau$  as:

$$\hat{\tau} = \tau + \delta_1 \frac{Cov(multiple\ birth, fertility\ treatment)}{Cov(multiple\ birth, family\ size)} \quad (5)$$

Equation (5) demonstrates that the bias of the IV estimate will depend on two elements: Firstly, the strength of the relationship between fertility treatments and the respective outcome ( $\delta_I$ ), i.e., how strongly the differences between mothers with and without fertility treatment affect the outcome of interest, and secondly, the importance of fertility treatments for the occurrence of multiple births - the covariance between multiple births and fertility treatments. This covariance is likely to be positive as the use of fertility treatments is consistently linked to multiple births in the medical literature (e.g., Callahan et al., 1994; Gleicher et al., 2000; Fauser, Devroey and Macklon, 2005). Indeed, as stated earlier, in our sample the likelihood of having multiple births is 1% for women without fertility treatment and 13% for women who had fertility treatment and 24% of all multiple births observed in the data are preceded by fertility treatments.

When an increasing number of women use fertility treatments, the second part of the bias term in (5) will become larger as  $Cov(multiple\ birth, fertility\ treatment)$  will increase. It is also possible that  $\delta_I$  will change as the composition of the group of women who undergo fertility treatment changes, with  $\delta_I$  being zero if either no or all multiple births are due to fertility treatments. Furthermore, it is not possible, *a priori*,

to sign  $\delta_I$ . For example, in labour supply regressions,  $\delta_I$  could be positive because fertility treatments are used by individuals with a higher propensity to work, or it could be negative as the use of fertility treatments will be correlated with a desire for children, which might in turn be correlated with fewer individuals choosing employment.

Faced with these problems there are two ways to block the backdoor path  $Y \leftarrow X \rightarrow \text{fertility treatment} \rightarrow \text{multiple birth} \rightarrow \text{family size} \rightarrow Y$  opened by the relationship between  $X$ , *fertility treatment* and *multiple birth*. Firstly, if we observe fertility treatment, as we do, then it is possible to condition on it directly. This closes the backdoor path and removes any association between the confounders in  $X$  and *multiple birth*. Secondly, if all elements in  $X$  were observed, one could condition on those directly, which would have an equivalent effect. A problem with this second strategy is that it is unlikely that all elements of  $X$  are observed in any given dataset. However, as the first option is only available when the use of fertility treatments is observed, conditioning on variables that may be part of  $X$  may be the only option when using datasets lacking this information. This strategy has its own risk as it may introduce further bias, rather than ameliorating the existing bias. Theoretically, it is only clear that conditioning on the full set of confounders in  $X$  would cause  $\delta_I$  to be zero and eliminate the bias. Conditioning on a subset of confounders can attenuate the problem if  $\delta_I$  shrinks towards zero as a result. However, it could also aggravate the problem. Consider a case where  $X$  consists of only two variables,  $A$  and  $B$ , whose effects cancel each other out, so that  $\delta_I$  would be zero without conditioning. Conditioning on either one of them in this case would cause  $\delta_I$  to be non-zero and would actually increase bias.

### *C. Time varying treatment effects*



A second issue with the use of twin-birth instruments concerns the timing of the twin births or (equivalently) the age of twins in the sample. There are three issues to consider: i) The impact of a twin birth on family size if people have the opportunity to adjust their fertility over time (the first stage), ii) the impact of the twin birth on labour supply, which might change over time as the twins age (the reduced form) and iii) the age distribution of twins in the respective population, which will determine the overall effect in an IV labour supply regression, as it determines the weights in the aggregation of individual-level effects to an overall effect.

Some of these issues have been considered previously in the literature, but the problems in their entirety, and their possible link to IVF, have not been fully discussed. Jacobsen et al. (1999) highlight the fact that many families will adjust their subsequent fertility decisions to compensate for the presence of twins. To illustrate this point, consider the first stage:

$$Family\ size_i = \pi + \gamma^* multiple\ birth_i + \mu_i, \quad (6)$$

The logic behind the instrument is that the birth of twins leads to a larger-than-planned family size. In other words, the instrument only works if families cannot (fully) adjust to the arrival of an additional child. An example where this condition would be fulfilled is a woman giving birth to twins at the last planned birth, i.e., a case where a woman who wanted one further last child receives two instead. However, it is important to note that there will be a substantial number of women for whom realised fertility in the long term is unaffected by twin births. Whenever a twin birth occurs at any birth before the last, it is, in principle, possible to adjust fertility over the following years. Say a woman always wanted two children. At her first planned pregnancy she gives birth to twins. This twin birth will have different effects in the short and the long term. In the short term, she has one more child than she

planned to have at this point in time. In the long term, however, she can decide not to have another child and can end up with her originally planned family size. This suggests that the first stage could be written as

$$\begin{aligned} \text{Family size}_{it} = & \pi + \gamma_t * \text{multiple birth}_{it} + \gamma_{t-1} * \text{multiple birth}_{it-1} + \gamma_{t-2} * \text{multiple birth}_{it-2} \\ & + \dots + \gamma_{t-k} * \text{multiple birth}_{it-k} + \mu_{it}, \end{aligned} \quad (7)$$

where we allow the effect a multiple birth to be different dependent on when it occurred in relation to the point in time family size is measured. Equation (7) highlights the fact that individuals can adjust their family size after a multiple birth. For households, we expect the impact of multiple births on family size to fall over time as found by Rosenzweig and Wolpin (1980), Bronars and Groggar (1994) and Jacobsen et al. (1999).

The exact value of  $\gamma_{tk}$  depends on the share of the multiple births being twins, triplets, quadruplets, etc. If there were only twins,  $\gamma_{tk}$  would be 1. As the vast majority of multiple births tend to be twins (96% in our sample), the estimate of  $\gamma_{tk}$  should be close to 1 directly after birth. A direct implication of this adjustment of fertility is that the composition of the complier group, i.e., those individuals who have a larger-than-planned family at each point in time, might change over time.

Similarly, we would expect the reduced form, i.e., the impact of a multiple birth on labour market outcomes, to weaken over time as children grow up, become more independent and finally leave the household. This suggests that the reduced form could be written as

$$\begin{aligned} Y_{it} = & \alpha + \lambda_t * \text{multiple birth}_{it} + \lambda_{t-1} * \text{multiple birth}_{it-1} + \lambda_{t-2} * \text{multiple birth}_{it-2} \\ & + \dots + \lambda_{t-k} * \text{multiple birth}_{it-k} + \varepsilon_{it}, \end{aligned} \quad (8)$$

This model captures the fact that the impact of the multiple-birth induced variation in family size may change over time, for example as the children become less dependent

on their mother as they grow up. As the second stage is simply the reduced form divided by the first stage,

$$\tau = \lambda/\gamma, \tag{9}$$

the estimated treatment effect would be expected to change with the time passed since birth as well.

In the following section we present tests for changes in the composition of the complier group and comparisons of first stages and labour supply estimates for a single birth cohort 1, 3 and 5 years after birth. Implications for estimates based on single cross-sections will be discussed in section 5.B after the presentation of results.

#### *D. Modelling*

We estimate and compare six models across our three samples collected at different intervals after birth. We have framed the discussion in this section in terms of a continuous outcome  $Y$  as most of the literature uses linear models. We have also estimated instrumental variable probits for binary outcomes, such as whether the individual is employed. The magnitude of the results is comparatively similar to the 2SLS results that we present. More importantly, the relative pattern of results across the different models, which matters for this paper, is practically identical. In other words, using an IV probit instead of 2SLS (unsurprisingly) does not help at all with an eventual bias caused by fertility treatments being unobserved.

The first model uses information that would be available in most datasets, meaning we instrument for family size using a dummy for whether the woman gave birth to twins or triplets ignoring the information on family size. The second includes a control for whether she also received fertility treatment. We also estimated models where we conditioned on the type of fertility treatment. These did not lead to any changes relative to a model only including a dummy for having received fertility

treatment. Estimates from the model controlling for fertility treatments are consistent since conditioning on fertility treatments is sufficient for the multiple-birth instrument to be valid. A comparison of these two models provides a picture of the size of the bias caused by unobserved fertility treatments. The third model considers the situation where information on fertility treatment is unavailable by conditioning on a set of variables that should be available in most datasets and that could plausibly be part of  $X$ . These include the education of the mother, whether she worked before the pregnancy, age at birth, ethnicity and marital status. Given the relative richness of information in the MCS we could condition on additional variables. However, we deliberately restrict our choice to variables that are realistically available to researchers trying to use the multiple-birth instrument with standard household data. A comparison of this model with the two previous models allows us to judge whether this conditioning strategy helps to attenuate any eventual bias. In a fourth model, we condition on both fertility treatments and the previously mentioned pre-pregnancy characteristics. Our discussion suggests that the only link between  $X$  and *multiple birth* arises due to fertility treatments. If this is indeed the case, conditioning on pre-pregnancy characteristics and fertility treatments should not lead to different results than conditioning on fertility treatments alone. Finally, we also evaluate whether the first and second stages for women with and without fertility treatments are different by estimating separate models for the two groups and comparing the results. All of these estimates include dummy variables for the current age of the mother in years to control for the earlier discussed age differences between mothers with single and multiple births.

We also test for differences in the characteristics of the compliers within a single birth cohort over time. Compliers are generally unobservable in the data,

however, there are ways to characterize them (Angrist and Pischke, 2009, pp. 166-172): For discrete characteristics  $x_i$ , we can describe the likelihood of a complier having that characteristic relative to the population by dividing the first stage for the sub-sample with  $x_i = 1$  by the overall first stage. The resulting complier-population ratios should be interpreted as relative likelihoods, i.e., a value of 2 indicates that compliers are twice as likely to have the respective characteristic than the general population. Values above 1 indicate that the characteristic is more common among the compliers than in the population and values below 1 indicate the opposite. All the characteristics we consider are based on pre-pregnancy characteristics, i.e., they are by construction unaffected by a later multiple or singleton birth. We repeat this exercise for all three sweeps of our data and compare results.

#### *E. Descriptive comparisons*

[TABLE 2 ABOUT HERE.]

Table 2 compares the pre-pregnancy characteristics of women based on sweep 1 of the MCS. There are a range of statistically significant and economically large differences. Women with fertility treatment are on average older (4 years older at the time of recorded birth), are more likely to have a (higher or first) degree and are less likely to have no qualifications. In terms of work and marital status, women with fertility treatment are 20 percentage points more likely to have worked before the pregnancy, are more likely to be married and are less likely to be single. Most women with fertility treatment are white, are 6 percentage points less likely to be non-white and have somewhat smaller families at sweep 1 (despite the higher likelihood of multiple births). Furthermore, those with fertility treatment are 13 percentage points more likely to have experienced complications during pregnancy and are 47 percentage points less likely to have an unplanned pregnancy. For most of these

factors it is easy to imagine a link with labour supply. Regressing a dummy for having received fertility treatment on these characteristics results in an  $R^2$  of about 0.04, suggesting that these are not the only characteristics in which women with and without fertility treatment differ.

[TABLE 3 ABOUT HERE.]

As stated before it should be possible to use multiple births as an instrument after conditioning on fertility treatments, as multiple births are probably still conditionally random. Table 3 provides some evidence on this conjecture. We compare the same characteristics as in Table 2 between women with singleton and multiple births conditional on having received fertility treatment. While there are still some significant differences between women with single and multiple births in each group, these are generally a lot smaller and often not statistically significant. These suggest that using multiple births as an instrument for family size might be possible as long as we are able to condition on having undergone fertility treatment.

#### **4 Female labour supply**

We begin by documenting differences in the outcomes between women with and without fertility treatments conditional on having had a single or multiple births.

[TABLE 4 ABOUT HERE.]

Table 4 shows these differences: In general, single-birth women with and without fertility treatment appear to differ in various dimensions: Women with fertility treatment are more likely to have a working partner in all sweeps and are also significantly more likely to be working in both sweeps 1 and 2. They are also more likely to use paid childcare. The differences in employment seem to largely disappear by sweep 3 when most, 99%, of the children in our data attend school. For those who

work, working hours do not appear to be too different. Women with multiple births in the two groups appear to be much more similar. While there are still differences in the probability of having a working partner in all sweeps, the gap in employment probabilities is smaller than among single-birth women and only significantly different from zero in sweep 1. These results suggest that there are some differences between the groups that are not related to variations in family size caused by multiple births. We now evaluate whether these also lead to differences in the first and second stages of standard labour supply regressions.

[TABLE 5 ABOUT HERE.]

Table 5 presents the first stage regression results. Consider first the models in columns (i) and (ii). The inclusion of a control for fertility treatment strengthens the estimated relationship between multiple births and family size: The coefficient on multiple births increases by about 20% in sweeps 1 and 2 and by about 25% in sweep 3. At the same time, the first stage F-value increases substantially. Conditioning on pre-pregnancy characteristics in column (iii) strengthens the first-stage relationship, but does very little to the first-stage coefficient on multiple births relative to column (i). The results from column (iv), where we condition on both fertility treatments and pre-pregnancy characteristics, are very similar to column (ii), but with a slightly higher F-value. The latter is simply the familiar result that IV estimates improve in precision after conditioning on other exogenous variables.

Comparing the first stages for women with and without fertility treatment in columns (v) and (vi) reveals that the instrument is a much better predictor of family size for women with fertility treatments with much higher first stage  $R^2$ -values and similar F-values despite a much smaller sample size.

Finally, the evidence in table 5 suggests that the time passed since the twin birth matters for the results: 1 year after the birth the impact of a multiple birth on family size (measured at the time of the respective survey) are substantially larger than in later sweeps. In fact, in the models that are likely to be unbiased, the impact on family size is slightly above 1. This is sensible given that women would not have had time to adjust their future fertility in response to the multiple births and that a twin birth would result in one extra child, while the few triplets in our data would result in 2 extra children. In later sweeps women have been able to make adjustments to their fertility enabling some of them to return to their planned family size. However, the instrument remains strong with a positive effect on family size, suggesting that a substantial share of mothers end up with more children than they originally intended.

[TABLE 6 ABOUT HERE.]

Table 6 presents results from the characterisation of compliers in each of the three samples, i.e., those mothers with a larger-than-planned family size due to having experienced a multiple birth. In sweep 1, compliers appear to be more likely to have had a surprising pregnancy and have either no or relatively high qualifications. Compliers are less likely to have medium qualifications such as O-levels/GCSEs and A-levels. The former are the first school-leaving qualification pupils can take, usually at the ages of 14 to 16. A-levels are further education qualifications taken at the age of 18 and are usually required for university admittance. Compliers are also less likely to have experienced problems during the pregnancy. In terms of ethnicity, marriage and employment before the pregnancy, compliers and non-compliers in the sample are quite similar.



Over time, we can see marked changes in the composition of the compliers: In sweep 3, compliers are more likely to have low and medium qualifications up to A-level, and increasingly less-likely, relative to sweep 1, to have a diploma or a degree. Compliers in sweep 3 are also much more likely to be non-white than the general population and have about the same share of people who experienced problems during the pregnancy. Furthermore, the proportion of compliers experiencing surprise pregnancies increases relative to sweep 1. Overall, the evidence suggests that the composition of compliers changes quite markedly with time passed since the respective multiple birth.

[TABLES 7, 8 AND 9 ABOUT HERE.]

By estimating our models for each sweep we can investigate how the results change across time with the changing complier groups. Tables 7 to 9 present evidence for sweeps 1, 2 and 3, respectively. The results in columns (i) and (ii) are generally similar. Having more children lowers the propensity to be working, in sweep 2 also that of self-employment, and increases the probability of staying at home and caring for the family. These effects also appear to be stronger when at least one of the children is young and decline as the child ages (across sweeps 1 to 3). There also does not appear to be any effect on the working hours for those who are working. The relatively similarity of the results in these two columns suggest that the bias from omitting fertility treatments might be negligible.

The results from column (iii) suggest that conditioning on pre-pregnancy characteristics also does not lead to substantial changes in results. However, there are several cases where the size of coefficients in column (iii) is different from those in both columns (i) and (ii). This finding highlights that conditioning on a subset of potential confounders might sometimes make matters worse. The results in column

(iv) generally suggest that adding pre-pregnancy characteristics does not change the results if we also condition on fertility treatments. This result again suggests that the only source of correlation between multiple births and mothers' characteristics arises because of fertility treatments.

The third thing to note is that in columns (v) and (vi) the magnitude of the effects seems fairly similar for women with and without fertility treatment. In sum, the results suggest that despite existing differences between women with and without fertility treatment, the bias in labour supply regressions relying on a multiple-birth instrument appears to be comparatively small.

Comparing results across the three sweeps suggest very different effects on labour supply. For sweep 1, the effect of the twin-birth induced variation in family size on female employment is strongly negative and both economically and statistically significant. We also see that the relationship between number of children and staying at home to care for family is positive and significant. Three years after the birth, in sweep 2, the effects are similar in magnitude, even though they have become weaker in terms of statistical significance. After 5 years, however, the picture changes substantially, point estimates are much closer to zero and are always insignificant.

## **5 Discussion and implications**

### **A. Unobserved fertility treatments**

Our results clearly suggest that omitting fertility treatments has very different effects on the first and second stages of the IV regression. First stages appear to be downward-biased due to the fact that fertility treatments are more likely to be taken by women with otherwise smaller families. Adjusting for fertility treatments strengthens the first stage. Interestingly, this bias does not appear to carry over to the second stages even though mothers with and without fertility treatments differ in a

range of characteristics that would make a second stage bias likely. Adjusting for a range of pre-pregnancy characteristics that one could realistically observe in most household data sets does not appear to help and in fact seems to increase bias slightly in some of our specifications. Overall, these estimates suggest that multiple births might still be a reasonable instrument for family size in labour supply regressions, even in countries and time periods where fertility treatments are common.

It is clear from our estimates that the main reason why we do not observe any bias is that the second stage coefficients in the two groups of mothers are very similar. What is less clear from our estimates is why this is the case given the observed differences between mothers with and without fertility treatment in characteristics such as education, pre-pregnancy work experience and health.

## **B. Time-varying effects**

The results in section 4 show that the effects of the multiple-birth induced variation in family size depend on the time passed since the multiple birth. First stages change as people adjust their fertility, leading to changes in the composition of the compliers, which, combined with the effects of children growing older and becoming more independent, leads to marked differences in the second stage coefficients.

One implication arising from this observation is that estimating and comparing labour supply regressions across different cross-sectional samples without accounting for the time passed since the multiple births might be problematic. If we estimate the first stage as in (6) as

$$Family\ size_i = \pi + \gamma^* multiple\ birth_i + \mu_i, \quad (10)$$

$\gamma$  is a weighted average of the  $\hat{\gamma}$  that would result from estimating equation (7).

Correspondingly, the reduced form coefficient  $\hat{\lambda}$  is a weighted average of the  $\hat{\lambda}_t$ s from equation (8). The weights in both cases depend on the age distribution of the

children resulting from multiple births. If the age distribution of these children was constant over time, comparisons between estimates based on different samples would not be problematic as the weighting of the first stage and reduced form coefficients would be identical in the different samples. However, if, as we observe, multiple births are increasing over time, then  $\gamma$  may be larger in later cohorts than earlier cohorts, not because the impact of multiple birth on family size at the individual level is changing, but because the number of younger twins is increasing in the population. Such differences mean that comparisons of results from different cross-sections would be affected by the distribution of twins.

(FIGURE 3 ABOUT HERE.)

Consider, for example, a case where a researcher has access to cross-sectional datasets, say a (hypothetical) census conducted in 2000 and 2010. Twins in a census can be identified as long as they live in their parents' house, for simplicity assume that this occurs up to the age of 20. The estimates based on the 2000 census would then effectively rely on twin births that occurred during the period 1980 to 2000. This situation is depicted in panel 3(a) of Figure 3, where the dashed lines mark this period. For the 2010 census, estimates would be based on twin births from 1990 to 2010. This is illustrated in panel 3(b). If the effects of twin births vary over time, either because effects genuinely vary with the time passed since birth or because the composition of compliers in each birth-cohort changes over time, the estimates in the first and second stages will depend partially on the distribution of twins across birth cohorts and time. If more twin births occurred relatively close to the census date, the estimated effects are likely to be dominated by the short-term effects, i.e., a combination of relatively fewer families being able to adjust their family size and relatively young children in the families experiencing a multiple birth. If a larger

proportion of the multiple births in the population occurred earlier, however, first stages are likely to be weaker as more families have had time to adjust their fertility. Similarly, as the children born in the multiple births would be older, the reduced form coefficients might also be closer to zero. As the second stage is simply the reduced form divided by the first stage, i.e.  $\tau = \lambda/\gamma$ , the estimated treatment effect in the latter case could be larger or smaller than in the first case.

When comparing cross-sectional results, such as in our 2000 and 2010 census example above, there are in principle several explanations for differences in the estimates. First, the effect of family size on female labour supply might have changed, be it because of changes on the individual level, such as attitudes, or be it because of changes to public policy, such as child care. A second explanation would be that the distribution of twin births over time (i.e., the age distribution of twins) in the two samples is different, leading to a different weighting of the time-varying effects of the twin-birth-induced fertility. A third possible explanation is a change in the composition of the compliers in both samples that is unrelated to the time passed since birth. Furthermore, if the frequency of multiple births in the population is related to IVF decisions, the endogeneity problem we discussed earlier might also be more or less severe in one of the two samples. While these arguments do not necessarily point towards a “bias” in the conventional definition, they are definitely another source of heterogeneity that hinders the comparison of results across papers using different samples.

## 6 Conclusion

This paper evaluated the rise of fertility treatments as a threat to the commonly used multiple-birth instrument for family size. Fertility treatments might threaten this

identification strategy as they are linked to the occurrence of multiple births as well as to a range of characteristics that might influence labour supply. Using the British Millennium Cohort Study, which allows us to distinguish between women that have and have not used fertility treatment, we investigate the impact of family size on labour supply outcomes with and without controlling for fertility treatment.

We find that there are indeed differences, both in pre-pregnancy characteristics and outcomes, between women that have and have not used fertility treatments. Conditional on having undergone fertility treatment, the birth of twins or triplets appears to be a random event. Omitting fertility treatment appears to bias first stages downward, weakening the first stage relationship. The bias in the second stages that arises from omitting fertility treatment controls appears to be comparatively small in magnitude and does not affect qualitative results. In all specifications, conditioning on a set of typically observed pre-pregnancy characteristics, rather than fertility treatment itself, does not appear to help very much and might in fact slightly increase bias.

We also find evidence that effects depend strongly on the time passed since the birth of the twins: First stages become weaker over time even though the instrument remains strong throughout. We also observe that the composition of compliers changes as individuals adjust their fertility over time. Second stages change considerably between regressions at 9 months and 3 and 5 years after the births, with point estimates getting closer to zero and increasingly becoming statistically insignificant. This pattern of results implies that one might get very different results from a dataset where most of the twin births occurred several years before the sampling period than from one where most twin births are relatively recent. It also

suggests that estimates from any cross-sectional dataset will always depend on the distribution of birth dates for the twins (or triplets) in the sample.

### References

Angrist, J. D. and Evans, W. N. (1998). Children and their parents' labor supply: Evidence from exogenous variation in family size. *American Economic Review* 88(3), 450-477.

Angrist, J. D. and Imbens, G. W. (1994). Identification and estimation of local average treatment effects. *Econometrica* 62(2), 467-475.

Angrist, J. D., Imbens, G. W. and Rubin, D. B. (1996). Identification of causal effects using instrumental variables. *Journal of the American Statistical Association* 91(434), 444-455.

Angrist, J. D. and Pischke, J (2009). *Mostly harmless econometrics – an empiricist's companion*. Princeton University Press.

Angrist, J., Lavy, V. and Schlosser, A. (2010). Multiple experiments for the causal link between the quantity and quality of children. *Journal of Labor Economics* 28(4), 773-823.

Black, S., Devereux, P. and Salvanes, K.G. (2005). The more the merrier? The effect of family composition on children's outcomes. *The Quarterly Journal of Economics* 120(2), 669-700.

Bronars, S. G. and Grogger, J. (1994). The economic consequences of unwed motherhood: Using twin births as a natural experiment. *American Economic Review* 84(5), 1141-56.

Callahan, T. L., Hall, J.E., Ettner, S.L., Christiansen, C.L., Greene, M.F. and Crowley, W.F. (1994). The economic impact of multiple-gestation pregnancies and the contribution of assisted-reproduction techniques on their incidence. *New England Journal of Medicine* 331(4), 244-249.

Dex, S. and Joshi, H. (2005) *Children of the 21st century: from birth to nine months*, Policy Press, Bristol, UK

Fauser, B. C., Devroey, P. and Macklon, N.S. (2005). Multiple birth resulting from ovarian stimulation for subfertility treatment. *The Lancet* 365(9473), 1807-1816.

Gleicher, N., Oleske, D.M., Tur-Kaspa, I., Vidali, A. and Vishvanath Karande (2000). Reducing the risk of high-order multiple pregnancy after ovarian stimulation with gonadotropins. *New England Journal of Medicine* 343(1), 2-7.

Hansen, K. and Joshi, H. (2007) *Millennium Cohort Study Second Survey: a user's guide to initial findings*, Institute of Education, London, UK

Human Fertilisation and Embryology Authority (2012) *Fertility Treatment in 2012: Trends and Figures*, HFEA, UK



Jacobsen, J. P., Wishart Pearce III, J. and Rosenbloom, J.L. (1999). The Effects of Childbearing on Married Women's Labor Supply and Earnings: Using Twin Births as a Natural Experiment. *Journal of Human Resources* 34(3), 449-474.

Morgan, S. L. and Winship, C. (2007). *Counterfactuals and causal inference*. Cambridge University Press, Cambridge.

Pearl, J. (2000). *Causality: Models, Reasoning, and Inference*. Cambridge University Press, Cambridge.

Qualifications and Curriculum Authority (2003). *Foundation Stage Profile Handbook*. London.

Rosenzweig, M. and Wolpin, K.I. (1980). Testing the quantity-quality fertility model: The use of twins as a natural experiment. *Econometrica* 48(1), 227-240.

Table 1: Descriptive statistics labour supply sample

Variable	Observations	Mean	Std.dev.	Min.	Max.
Twin birth	18340	0.01	0.11	0	1
Triplet birth	18340	0.001	0.02	0	1
Multiple birth	18340	0.01	0.12	0	1
Had fertility treatment	18340	0.03	0.16	0	1
Pregnancy was surprising	18340	0.46	0.50	0	1
No qualification	18340	0.20	0.40	0	1
Qualification up to O-level/GCSE or equivalent	18340	0.34	0.47	0	1
A-level	18340	0.09	0.29	0	1
Higher education diploma	18340	0.08	0.28	0	1
First degree	18340	0.12	0.33	0	1
Higher degree (Master, PhD)	18340	0.03	0.18	0	1
Age at birth	18340	28.3	5.95	14	51
Had job before pregnancy	18340	0.02	0.15	0	1
Non-white ethnicity	18340	0.160	0.37	0	1
Married (1 <sup>st</sup> marriage)	18340	0.56	0.50	0	1
Remarried (2 <sup>nd</sup> or higher marriage)	18340	0.04	0.20	0	1
Single	18340	0.34	0.47	0	1
Divorced or separated	18340	0.07	0.25	0	1
Illness or problems during pregnancy	18340	0.38	0.48	0	1
Fertility and outcomes at time of sweep 1 interview (within 1 year of birth)					
Number of children	18340	2.0	1.09	1	10
Age	18340	29.1	5.95	14	52
Employed	18340	0.40	0.49	0	1
On maternity leave	18340	0.02	0.13	0	1
Self-employed	18340	0.03	0.16	0	1
Student	18340	0.01	0.09	0	1
At home to care for family	18340	0.54	0.50	0	1
Weekly working hours (includes 0)	18340	11.7	14.66	0	86
Weekly working hours (excludes 0)	8669	24.8	11.35	1	86
Has working partner	18340	0.72	0.45	0	1
Fertility and outcomes at time of sweep 2 interview (3 years after birth )					
Number of children	14460	2.2	1.08	1	13
Age	14460	31.9	5.85	17	54
Employed	14460	0.48	0.50	0	1
Self-employed	14460	0.01	0.09	0	1
Student	14460	0.01	0.11	0	1
At home to care for family	14460	0.44	0.50	0	1
Weekly working hours (includes 0)	14460	12.4	14.38	0	114
Weekly working hours (excludes 0)	7558	23.8	11.19	1	114
Has working partner	14460	0.75	0.43	0	1
Uses childcare by conducted by relatives/friends	14460	0.28	0.45	0	1
Uses paid childcare	14460	0.13	0.33	0	1
Fertility and outcomes at time of sweep 3 interview (5 years after birth)					
Number of children	12581	2.4	1.06	1	13
Age	12581	34.1	5.81	18	58
Employed	12581	0.53	0.50	0	1
Self-employed	12581	0.01	0.11	0	1
Student	12581	0.01	0.11	0	1
At home to care for family	12581	0.37	0.48	0	1
Weekly working hours (includes 0)	12581	14.0	14.46	0	100
Weekly working hours (excludes 0)	7390	23.8	11.10	0	100
Has working partner	12581	0.75	0.43	0	1
Child attends school	12581	0.99	0.11	0	1

Table 2: Comparison of pre-pregnancy characteristics of women with and without fertility-treatment

Variable	<u>Without fertility treatment</u>		<u>With fertility treatment</u>		P-Value means different <sup>a</sup>
	Mean	Std.dev.	Mean	Std.dev.	
Twin birth	0.01	0.10	0.12	0.32	0.00
Triplet birth	0.00	0.02	0.01	0.10	0.03
Multiple birth	0.01	0.10	0.13	0.33	0.00
Pregnancy was surprising	0.47	0.50	0.00	0.00	0.00
Birth weight 1 <sup>st</sup> child (kg)	3.36	0.57	3.19	0.65	0.00
Number of children at sweep 1 interview	1.96	1.09	1.54	0.75	0.00
No qualification	0.20	0.40	0.12	0.32	0.00
Qualification up to O-level/GCSE or equivalent	0.34	0.47	0.33	0.47	0.98
A-level	0.09	0.29	0.11	0.31	0.33
Higher education diploma	0.08	0.28	0.09	0.29	0.44
First degree	0.12	0.33	0.18	0.38	0.00
Higher degree (Master, PhD)	0.03	0.18	0.08	0.27	0.00
Age at birth	28.22	5.94	32.29	4.94	0.00
Had job before pregnancy	0.62	0.49	0.81	0.39	0.00
Non-white ethnicity	0.16	0.37	0.10	0.30	0.00
Married (1 <sup>st</sup> marriage)	0.55	0.50	0.76	0.43	0.00
Remarried (2 <sup>nd</sup> or higher marriage)	0.04	0.20	0.06	0.25	0.03
Single	0.34	0.47	0.12	0.32	0.00
Divorced or separated	0.07	0.25	0.06	0.24	0.40
Illness or problems during pregnancy	0.37	0.48	0.50	0.50	0.00
Observations	17,862		478		

<sup>a</sup> Based on two sample t-test with unequal variances.

Table 3: Comparison of pre-pregnancy characteristics of women with single and multiple births by fertility treatment

	Women without fertility treatment					Women with fertility treatment				
	Single birth		Multiple birth		P-Value means different	Single birth		Multiple birth		P-value means different
	Mean	Std.dev.	Mean	Std.dev.		Mean	Std.dev.	Mean	Std.dev.	
Pregnancy was surprising	0.47	0.50	0.49	0.50	0.58	0.00	0.00	0.00	0.00	
Birth weight 1 <sup>st</sup> child (kg)	3.37	0.56	2.44	0.52	0.00	3.30	0.58	2.42	0.59	0.00
No qualification	0.20	0.40	0.23	0.42	0.23	0.12	0.33	0.10	0.30	0.57
Qualification up to O-level/GCSE or equivalent	0.34	0.47	0.31	0.46	0.47	0.33	0.47	0.38	0.49	0.470
A-level	0.09	0.29	0.06	0.24	0.08	0.10	0.30	0.16	0.37	0.19
Higher education diploma	0.08	0.28	0.14	0.35	0.03	0.10	0.30	0.07	0.25	0.35
First degree	0.12	0.33	0.13	0.34	0.622	0.18	0.38	0.16	0.37	0.79
Higher degree (Master, PhD)	0.03	0.18	0.03	0.16	0.58	0.08	0.27	0.07	0.25	0.65
Age at birth	28.20	5.94	30.12	5.70	0.00	32.21	4.92	32.87	5.09	0.34
Had job before pregnancy	0.62	0.49	0.64	0.48	0.58	0.81	0.39	0.80	0.40	0.86
Non-white ethnicity	0.16	0.37	0.12	0.32	0.08	0.11	0.31	0.07	0.25	0.26
Married (1 <sup>st</sup> marriage)	0.55	0.50	0.59	0.49	0.25	0.75	0.44	0.85	0.36	0.04
Remarried (2 <sup>nd</sup> or higher marriage)	0.04	0.20	0.06	0.24	0.20	0.07	0.25	0.03	0.18	0.16
Single	0.34	0.47	0.28	0.45	0.09	0.12	0.33	0.05	0.22	0.02
Divorced or separated	0.07	0.25	0.06	0.24	0.75	0.06	0.23	0.07	0.25	0.81
Illness or problems during pregnancy	0.37	0.48	0.46	0.50	0.02	0.49	0.50	0.54	0.50	0.48
Observations	17,669		193			417		61		

Table 4: Comparisons of outcomes for women with and without fertility treatment with same number of children born

	Single births					Multiple births				
	No FT		FT		P-value means different	No FT		FT		P-value means different
	Mean	Std.dev.	Mean	Std.dev.		Mean	Std.dev.	Mean	Std.dev.	
Sweep I outcomes										
Employed	0.40	0.49	0.52	0.50	0.00	0.31	0.46	0.43	0.50	0.11
On maternity leave	0.03	0.16	0.06	0.23	0.01	0.02	0.14	0.03	0.18	0.63
Self-employed	0.02	0.13	0.04	0.19	0.02	0.04	0.19	0.08	0.28	0.23
Student	0.01	0.09	0.00	0.07	0.27	0.01	0.07	0.02	0.13	0.52
At home to care for family	0.54	0.50	0.38	0.49	0.00	0.63	0.48	0.44	0.50	0.01
Weekly working hours (includes 0)	11.63	14.61	16.78	15.62	0.00	9.98	14.49	15.36	15.18	0.02
Weekly working hours (excludes 0)	24.81	11.33	25.82	11.90	0.17	24.39	12.69	26.77	9.61	0.27
Has working partner	0.72	0.45	0.90	0.29	0.00	0.74	0.44	0.89	0.32	0.00
Observations	17,669		417			193		61		
Sweep II outcomes										
Employed	0.47	0.50	0.59	0.49	0.00	0.46	0.50	0.43	0.50	0.68
Self-employed	0.01	0.09	0.00	0.05	0.08	0.00	0.00	0.00	0.00	n/a
Student	0.01	0.11	0.01	0.08	0.12	0.01	0.12	0.06	0.24	0.20
At home to care for family	0.44	0.50	0.31	0.47	0.00	0.46	0.50	0.41	0.50	0.52
Weekly working hours (includes 0)	12.38	14.36	15.45	15.01	0.00	11.73	14.35	12.41	14.23	0.77
Weekly working hours (excludes 0)	23.83	11.14	23.34	12.49	0.56	23.14	11.89	24.35	10.13	0.62
Has working partner	0.75	0.43	0.92	0.28	0.00	0.77	0.42	0.82	0.39	0.45
Uses childcare by conducted by relatives/friends	0.28	0.45	0.29	0.46	0.72	0.23	0.42	0.16	0.37	0.27
Uses paid childcare	0.12	0.33	0.22	0.42	0.00	0.08	0.27	0.18	0.39	0.09
Observations	13,942		343			142		51		
Sweep III outcomes										
Employed	0.53	0.50	0.57	0.50	0.13	0.54	0.50	0.52	0.50	0.81
Self-employed	0.01	0.11	0.00	0.06	0.02	0.01	0.09	0.00	0.00	0.31
Student	0.01	0.11	0.01	0.08	0.22	0.01	0.09	0.02	0.14	0.57
At home to care for family	0.37	0.48	0.30	0.46	0.01	0.34	0.48	0.31	0.47	0.69
Weekly working hours (includes 0)	13.92	14.48	15.32	13.71	0.08	13.68	14.32	15.48	14.16	0.46
Weekly working hours (excludes 0)	23.80	11.11	22.53	10.67	0.10	22.86	11.47	23.97	10.20	0.63
Has working partner	0.75	0.43	0.89	0.32	0.00	0.80	0.41	0.88	0.33	0.19
Child attends school	0.99	0.11	0.99	0.11	0.85	0.98	0.13	0.96	0.20	0.42
Observations	12,111		300			122		48		

Table 5: First stage results: Effect of a multiple birth on family size with different sets of controls

	(i) All women	(ii) All women, controls for fertility treatment	(iii) All women, controls for pre-pregnancy characteristics	(iv) All women, controls for pre- pregnancy characteristics & fertility treatment	(v) Only women with fertility treatment	(vi) Only women without fertility treatment
Sweep I						
Multiple birth (1 = yes)	0.88*** (0.07)	1.04*** (0.07)	0.89*** (0.06)	1.02*** (0.06)	1.10*** (0.09)	1.03*** (0.09)
Fertility treatment (1 = yes)		-0.78*** (0.03)		-0.65*** (0.03)		
R <sup>2</sup>	0.01	0.02	0.23	0.24	0.25	0.01
Kleinbergen- Paap F-stat	149.1	215.7	194.5	265.3	136.5	141.21
Observations	18,340	18,340	18,340	18,340	478	17,862
Sweep II						
Multiple birth (1 = yes)	0.69*** (0.08)	0.83*** (0.08)	0.69*** (0.07)	0.81*** (0.07)	1.01*** (0.12)	0.79*** (0.09)
Fertility treatment (1 = yes)		-0.64*** (0.04)		-0.54*** (0.04)		
R <sup>2</sup>	0.01	0.02	0.19	0.20	0.18	0.01
Kleinbergen- Paap F-stat	78.9	116.3	107.7	147.7	69.1	71.5
Observations	14,460	14,460	14,460	14,460	394	14,066
Sweep III						
Multiple birth (1 = yes)	0.57*** (0.08)	0.72*** (0.08)	0.61*** (0.07)	0.73*** (0.07)	0.90*** (0.13)	0.67*** (0.10)
Fertility treatment (1 = yes)		-0.58*** (0.05)		-0.49*** (0.05)		
R <sup>2</sup>	0.00	0.01	0.16	0.16	0.14	0.00
Kleinbergen- Paap F-stat	48.0	73.1	68.0	95.4	50.5	41.6
Observations	12,581	12,581	12,581	12,581	348	12,233

Coefficient, robust standard errors in parentheses. \*/\*\*/\*\* denote statistical significance on the 10%, 5% and 1% level respectively. All estimates include age in years as dummies. Column (iii) also contains dummies for various completed qualifications, age at birth, a dummy for having worked before the pregnancy, a dummy for non-white ethnicity and dummy variables for marital status.

Table 6: Analysis of compliers characteristics, first stage coefficients for subsamples and relative frequency of compliers

Characteristic	Sweep 1		Sweep 2		Sweep 3	
	First stage	Relative frequency compliers	First stage	Relative frequency compliers	First stage	Relative frequency compliers
Full sample	0.88*** (0.07)		0.69*** (0.08)		0.57*** (0.08)	
Pregnancy was surprising	1.01*** (0.14)	1.15	0.98*** (0.16)	1.43	0.93*** (0.19)	1.62
No qualification	0.98*** (0.22)	1.12	0.80*** (0.26)	1.17	0.87*** (0.32)	1.52
Highest qualification O-level or equivalent	0.78*** (0.10)	0.88	0.708*** (0.12)	1.03	0.61*** (0.13)	1.07
Highest qualification A-level	0.74*** (0.14)	0.83	0.61*** (0.16)	0.90	0.58*** (0.18)	1.01
Highest qualification diploma	1.01*** (0.19)	1.15	0.66*** (0.19)	0.96	0.57*** (0.18)	1.00
Highest qualification degree	0.73*** (0.11)	0.82	0.41*** (0.12)	0.60	0.37*** (0.12)	0.65
Highest qualification higher degree (Master and PhD)	1.02*** (0.17)	1.15	0.68*** (0.19)	0.99	0.61*** (0.20)	1.06
Had job before pregnancy	0.89*** (0.07)	1.01	0.68*** (0.08)	0.99	0.57*** (0.08)	1.00
Non-white	0.89*** (0.24)	1.01	0.78*** (0.22)	1.14	0.71** (0.28)	1.23
Married	0.89*** (0.086)	1.01	0.63*** (0.09)	0.93	0.56*** (0.10)	0.97
Illness or problems during pregnancy	0.70*** (0.09)	0.791	0.59*** (0.10)	0.87	0.58*** (0.11)	1.02

Each cell is from a different regression. Displayed is the coefficient of “multiple birth” – the first-stage variable of interest – with robust standard errors in parentheses. \*/\*\*/\*\* denote statistical significance on the 10%, 5% and 1% level respectively. All estimates include age in years as dummies.

Table 7: Outcomes Sweep I interview (within 1 year of birth), second stage coefficients

Outcome	(i) All women	(ii) All women, controls for fertility treatment	(iii) All women, controls for pre-pregnancy characteristics	(iv) All women, controls for pre-pregnancy characteristics & fertility treatment	(v) Only women with fertility treatment	(vi) Only women without fertility treatment
Employed (1 = yes)	-0.11*** (0.03)	-0.11*** (0.03)	-0.12*** (0.03)	-0.11*** (0.03)	-0.10 (0.06)	-0.11*** (0.03)
Self-employed (1= yes)	-0.01 (0.01)	-0.01 (0.01)	-0.01 (0.01)	-0.01 (0.01)	-0.02 (0.02)	-0.01 (0.01)
On maternity/parental leave (1 = yes)	0.03* (0.02)	0.02 (0.01)	0.03* (0.01)	0.02 (0.01)	0.03 (0.03)	0.02 (0.01)
Fulltime student (1 = yes)	0.00 (0.00)	0.00 (0.00)	0.00 (0.01)	0.00 (0.00)	0.01 (0.01)	-0.00 (0.01)
At home and caring for family (1 = yes)	0.09*** (0.03)	0.10*** (0.03)	0.11*** (0.03)	0.10*** (0.03)	0.80 (0.06)	0.11*** (0.03)
Weekly working hours (includes 0 for those not working)	-2.00** (1.00)	-2.34*** (0.85)	-2.30*** (0.89)	-2.15*** (0.78)	-1.92 (1.83)	-2.45** (0.96)
Weekly working hours (excludes those not working)	-0.02 (1.28)	-0.19 (1.12)	0.11 (1.27)	-0.05 (1.11)	0.39 (1.70)	-0.61 (1.39)
Has a working partner (1= yes)	-0.01 (0.03)	-0.03 (0.02)	-0.03 (0.02)	-0.03 (0.02)	-0.03 (0.04)	-0.02 (0.03)
Observations (all but second working hours regression)	18,340	18,340	18,340	18,340	478	17,862
Observations (second working hours regression)	8669	8669	8669	8669	306	8363

Each cell is from a different regression. Displayed is the coefficient of “number of children” – the second-stage variable of interest – with robust standard errors in parentheses. \*/\*\*/\*\* denote statistical significance on the 10%, 5% and 1% level respectively. All estimates include age in years as dummies. Column (ii) additionally contains a dummy for having received fertility-treatment. Column (iii) also contains dummies for various completed qualifications, age at birth, a dummy for having worked before the pregnancy, a dummy for non-white ethnicity and dummy variables for marital status.



Table 8: Outcomes Sweep II interview (3 years after birth), second stage coefficients

Outcome	(i) All women	(ii) All women, controls for fertility treatment	(iii) All women, controls for pre-pregnancy characteristics	(iv) All women, controls for pre- pregnancy characteristics & fertility treatment	(v) Only women with fertility treatment	(vi) Only women without fertility treatment
Employed (1 = yes)	-0.08 (0.05)	-0.08* (0.04)	-0.10** (0.05)	-0.08* (0.04)	-0.19*** (0.07)	-0.04 (0.05)
Self-employed (1 = yes)	- 0.01*** (0.00)	-0.01*** (0.00)	-0.01*** (0.00)	-0.01*** (0.00)	-0.00 (0.00)	-0.01*** (0.00)
Fulltime student (1 = yes)	0.03 (0.02)	0.02 (0.01)	0.03 (0.02)	0.02 (0.01)	0.06* (0.03)	0.01 (0.01)
At home and caring for family (1 = yes)	0.07 (0.05)	0.08* (0.04)	0.09* (0.05)	0.08* (0.04)	0.12 (0.07)	0.07 (0.05)
Weekly working hours (includes 0 for those not working)	-2.47* (1.47)	-2.42** (1.23)	-2.70** (1.35)	-2.15* (1.15)	-4.00** (2.04)	-1.83 (1.49)
Weekly working hours (excludes those not working)	-0.81 (1.61)	-0.55 (1.43)	-0.47 (1.59)	-0.24 (1.41)	-1.08 (2.22)	-0.91 (1.76)
Has a working partner (1 = yes)	-0.01 (0.04)	-0.04 (0.04)	-0.04 (0.04)	-0.04 (0.03)	-0.08 (0.06)	-0.01 (0.04)
Observations (all except below)	14,460	14,460	14,460	14,460	394	14,066
Observations (second working hours regression)	7558	7558	7558	7558	253	7305

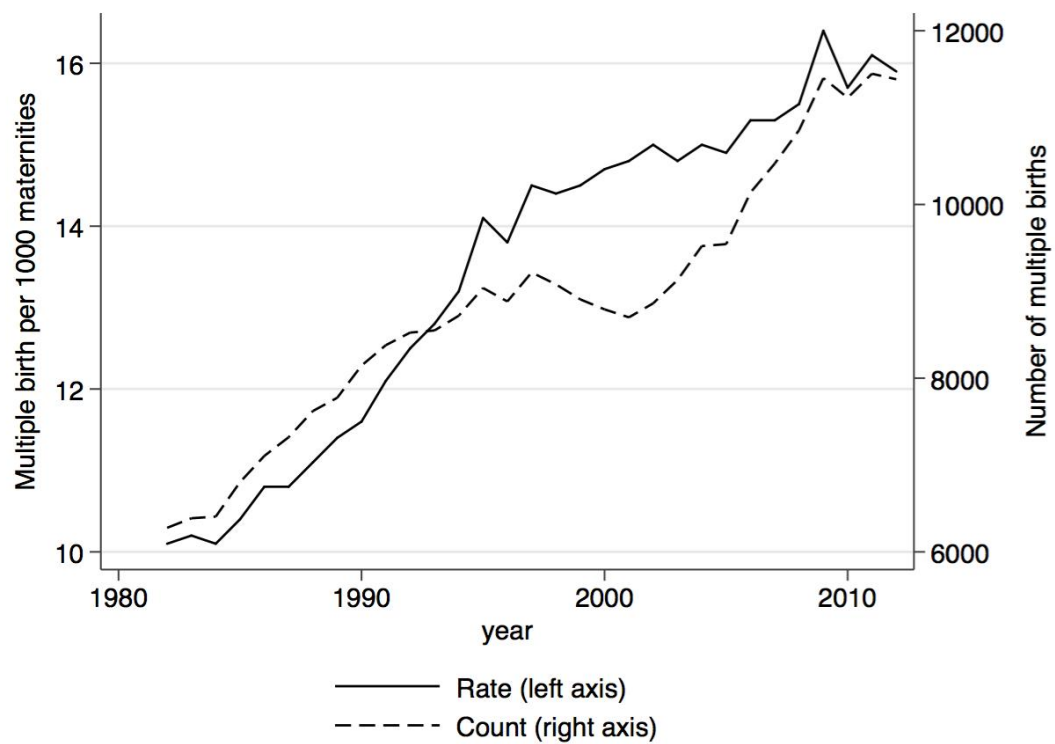
Each cell is from a different regression. Displayed is the coefficient of “number of children” – the second-stage variable of interest – with robust standard errors in parentheses. \*/\*\*/\*\* denote statistical significance on the 10%, 5% and 1% level respectively. All estimates include age in years as dummies. Column (ii) additionally contains a dummy for having received fertility-treatment. Column (iii) also contains dummies for various completed qualifications, age at birth, a dummy for having worked before the pregnancy, a dummy for non-white ethnicity and dummy variables for marital status.

Table 9: Outcomes Sweep III interview (5 years after birth), second stage coefficients

Outcome	(i) All women	(ii) All women, controls for fertility treatment	(iii) All women, controls for pre- pregnancy characteristics	(iv) All women, controls for pre- pregnancy characteristics & fertility treatment	(v) Only women with fertility treatment	(vi) Only women without fertility treatment
Employed (1 = yes)	-0.04 (0.07)	-0.03 (0.05)	-0.08 (0.06)	-0.04 (0.05)	-0.08 (0.09)	-0.01 (0.07)
Self-employed (1 = yes)	-0.01 (0.01)	-0.00 (0.01)	-0.00 (0.01)	-0.00 (0.01)	-0.00 (0.00)	-0.00 (0.01)
Fulltime student (1 = yes)	0.00 (0.02)	0.00 (0.01)	0.00 (0.01)	0.00 (0.01)	0.03 (0.03)	-0.00 (0.01)
At home and caring for family (1 = yes)	0.00 (0.06)	0.01 (0.05)	0.04 (0.06)	0.02 (0.05)	0.03 (0.08)	-0.00 (0.06)
Weekly working hours (includes 0 for those not working)	-1.59 (1.88)	-1.28 (1.53)	-2.37 (1.69)	-1.42 (1.42)	-0.55 (2.45)	-1.62 (1.91)
Weekly working hours (excludes those not working)	-1.27 (1.72)	-0.65 (1.46)	-1.00 (1.70)	-0.39 (1.44)	0.80 (2.34)	-1.42 (1.80)
Has a working partner (1= yes)	0.05 (0.05)	0.01 (0.04)	0.01 (0.05)	-0.00 (0.04)	-0.03 (0.06)	0.03 (0.05)
Observations (all but second working hours regression)	12,581	12,581	12,581	12,581	348	12,233
Observations (second working hours regression)	7390	7390	7390	7390	235	7155

Each cell is from a different regression. Displayed is the coefficient of “number of children” – the second-stage variable of interest – with robust standard errors in parentheses. \*/\*\*/\* denote statistical significance on the 10%, 5% and 1% level respectively. All estimates include age in years as dummies. Column (ii) additionally contains a dummy for having received fertility-treatment. Column (iii) also contains dummies for various completed qualifications, age at birth, a dummy for having worked before the pregnancy, a dummy for non-white ethnicity and dummy variables for marital status.

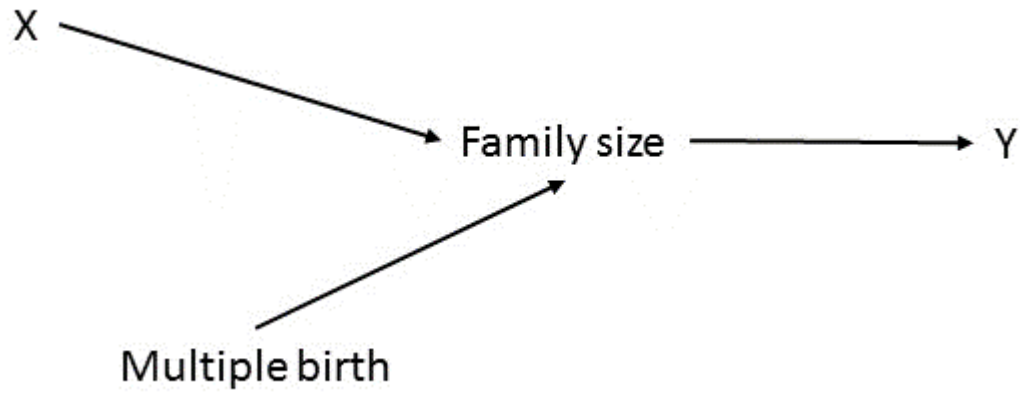
Figure 1: Multiple birth over time, UK, 1982 to 2012



Notes: Data is from the *Characteristics of Birth 2* series of the Office for National Statistics. We begin the series in 1982 as data from 1981 is missing due to a registrars' strike.

Figure 2: Causal diagram for the multiple-birth instrument with and without fertility treatments

Panel (a): The twin births instrument without fertility treatments



Panel (b): The twin births instrument with fertility treatments

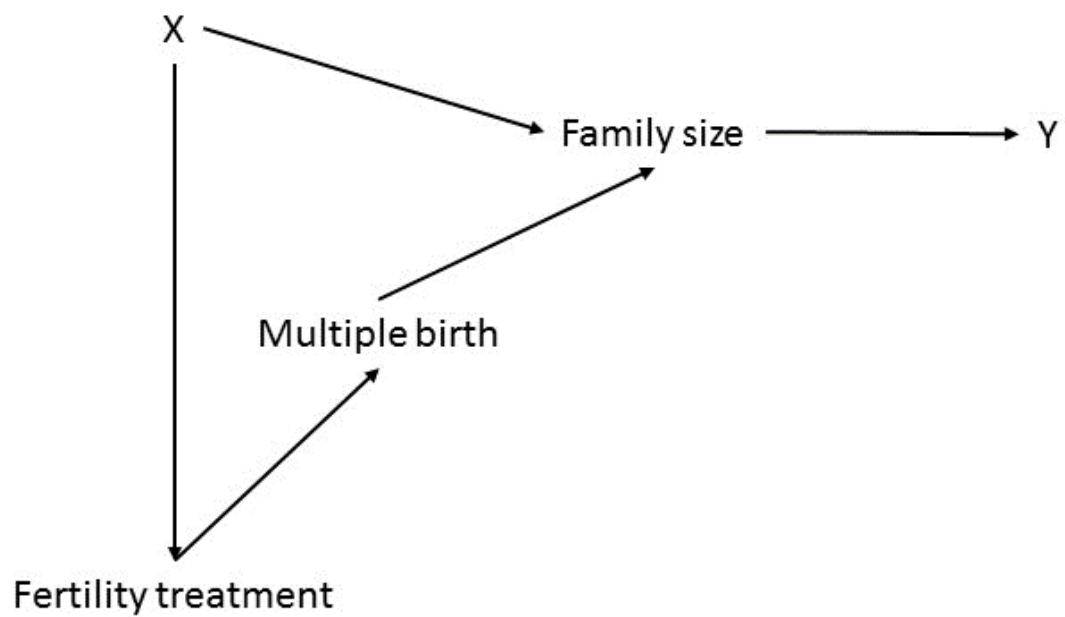
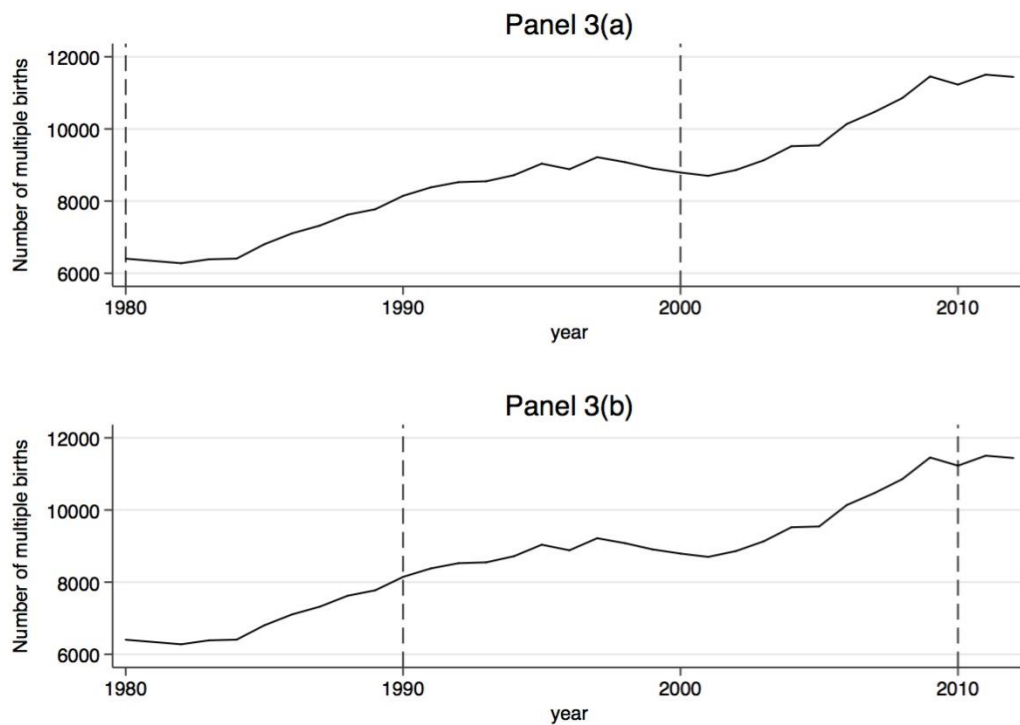


Figure 3: The distribution of multiple births over time and cross-sections drawn at various points



Notes: Data is from the *Characteristics of Birth 2* series of the Office for National Statistics. Data for 1981, which missing due to a registrars' strike, is linearly extrapolated between 1980 and 1982 for the sake of the example.